# Exact Analysis of Pattern Matching Algorithms with Probabilistic Arithmetic Automata

Tobias Marschall and Sven Rahmann

Bioinformatics for High-Throughput Technologies,
Computer Science XI, TU Dortmund, Germany
`firstname.lastname@tu-dortmund.de`

We propose a framework for the exact probabilistic analysis of window-based pattern matching algorithms, such as Boyer-Moore, Horspool, Backward DAWG Matching, Backward Oracle Matching, and more. In particular, we show how to efficiently obtain the distribution of such an algorithm's running time cost for any given pattern in a random text model, which can be quite general, from simple uniform models to higher-order Markov models or hidden Markov models (HMMs). Furthermore, we provide a technique to compute the exact distribution of *differences* in running time cost of two algorithms. In contrast to previous work, our approach is neither limited to simple text models, nor to asymptotic statements, nor to moment computations such as expectation and variance. Methodically, we use extensions of finite automata which we call *deterministic arithmetic automata* (DAAs) and *probabilistic arithmetic automata* (PAAs) [13]. To our knowledge, this is the first time that substring- or suffix-based pattern matching algorithms are analyzed exactly. Experimentally, we compare Horspool's algorithm, Backward DAWG Matching, and Backward Oracle Matching on prototypical patterns of short length and provide statistics on the size of minimal DAAs for these computations.

## 1 Introduction

The basic pattern matching problem is to find all occurrences of a *pattern* string in a (long) *text* string, with few character accesses. Let $\ell$ be the text length and $m$ be the pattern length. The well-known Knuth-Morris-Pratt algorithm [9] reads each text

character exactly once from left to right and hence needs exactly $\ell$ character accesses for any text of length $\ell$, after preprocessing the pattern in $\Theta(m)$ time. In contrast, the Boyer-Moore [4], Horspool [8], Sunday [19], Backward DAWG Matching (BDM, [5]) and Backward Oracle Matching (BOM, [1]) algorithms move a length-$m$ search window across the text and first compare its *last* character to the last character of the pattern. This often allows to move the search window by more than one position (at best, by $m$ positions if the last window character does not occur in the pattern at all), for a best case of $n/m$, but a worst case of $mn$ character accesses. The worst case can often be improved to $\Theta(m+n)$, but this makes the code more complicated and seldom provides a speed-up in practice. For practical pattern matching applications, the most important algorithms are Horspool, BDM (often implemented as Backward Nondeterministic DAWG Matching, BNDM, via a non-deterministic automaton that is simulated in a bit-parallel fashion), and BOM, depending on alphabet size, text length and pattern length; see [16, Section 2.5] for an experimental map.

A question that has apparently so far not been investigated is about the exact probability distribution of the number of required character accesses $X_\ell^p$ when matching a given pattern $p$ against a random text of finite length $\ell$ (non-asymptotic case), even though related questions have been answered in the literature. For example, [2, 3] analyze the expected value of $X_\ell^p$ for the Horspool algorithm. In [12] it is further shown that for the Horspool algorithm, $X_\ell^p$ is asymptotically normally distributed for i.i.d. texts, and [18] extends this result to Markovian text models. In [20], a method to compute mean and variance of these distributions is given.

In contrast to these results that are special to the Horspool algorithm, we use a general framework called *probabilistic arithmetic automata* (PAAs), introduced at CPM'08 [13], to compute the exact distribution of $X_\ell^p$ for any window-based pattern matching algorithm. In [13], PAAs where introduced in order to compute the distribution of occurrence counts of patterns; the fact that they can also be used to analyze pattern matching algorithms further highlights their utility. In our framework, the random text model can be quite general, from simple i.i.d. uniform models to high-order Markov models or HMMs. The approach is applied exemplarily to the following pattern matching algorithms in the non-asymptotic regime (short patterns, medium-length texts): Horspool, B(N)DM, BOM. We do not treat BDM and BNDM separately as, in terms of text character accesses, they are indistinguishable (see Section 2.2).

This paper is organized as follows. In the next section, we introduce notation and give a brief review of the Horspool, B(N)DM and BOM algorithms. In Section 3, we define *deterministic arithmetic automata* (DAAs). In Section 4, we present a simple general DAA construction for the analysis of window-based pattern matching algorithms. We also show that the state space of the DAA can be considerably reduced by adapting DFA minimization to DAAs. In Section 5, we summarize the PAA framework with its generic algorithms, define finite-memory text models and connect DAAs to PAAs. This yields, for each given pattern $p$, algorithm, and random text model, a PAA that computes the distribution of $X_\ell^p$ for any finite text length $\ell$. Section 6 introduces *difference DAAs* by a product construction that allows to compare two algorithms on a given pattern. Exemplary results on the comparison of several algorithms can be found in Section 7.

There, we also provide statistics on automata sizes for different algorithms and pattern lengths. Section 8 contains a concluding discussion.

An extended abstract of this work has been presented at LATA'10 [15] with two alternative DAA constructions. In contrast to that version, the DAA construction in the present paper can be seen as a combination of both of those, and is much simpler. Additionally, the DAA minimization introduced in the present paper allows the analysis of much longer patterns in practice. While [15] was focused on Horspool's and Sunday's algorithms, here, we give a general construction scheme applicable to any window-based pattern matching algorithm and discuss the most relevant algorithms, namely Horspool, BOM, and B(N)DM, as examples.

## 2 Algorithms

Both pattern and text are over a finite alphabet $\Sigma$. Indexing generally starts at zero. The pattern $p = p[0] \ldots p[m-1]$ is of length $m$; a (concrete) text $s$ is of length $\ell$. By $\overleftarrow{p}$, we denote the reverse pattern $p[m-1] \ldots p[0]$.

In the following, we summarize the Horspool, B(N)DM and BOM algorithms; algorithmic details can be found in [16, Chapter 2].

We do not discuss the Knuth-Morris-Pratt algorithm because its number of text character accesses is constant: Each character of the text is looked at exactly once. Therefore, $\mathcal{L}(X_\ell^p)$ is the Dirac distribution on $\ell$, i.e., $\mathbb{P}(X_\ell^p = \ell) = 1$.

We also do not discuss the Boyer-Moore algorithm, since it is never the best one in practice because of its complicated code to achieve optimal asymptotic running time. In contrast to our earlier paper [15], we also skip the Sunday algorithm, as it is almost always inferior to Horspool's. Instead, we focus on those algorithms that are fastest in practice according to [16, Fig. 2.22].

The Horspool, B(N)DM and BOM algorithms have the following properties in common: They maintain a search window $w$ of length $m = |p|$ that initially starts at position 0 in the text $s$, such that its rightmost character is at position $t = m - 1$. The right window position $t$ grows in the course of the algorithm; we always have $w = s[(t - m + 1) \ldots t]$. The algorithms look at the characters in each window from right to left, and thus compare the reversed window with the reversed pattern $\overleftarrow{p}$. (For Horspool, variants with different comparison orders are possible, but the rightmost character is always compared first.)

The two properties of an algorithm that influence our analysis are the following: For a pattern $p \in \Sigma^m$, each window $w \in \Sigma^m$ determines

1. its cost $\xi^p(w)$, e.g., the number of text character accesses required to analyze this window,

2. its shift $shift^p(w)$, which is the number of characters the window is advanced after it has been examined.

3

## 2.1 Horspool

First, the rightmost characters of window and pattern are compared; that means, $a := w[m-1] = s[t]$ is compared with $p[m-1]$. If they match, the remaining $m-1$ characters are compared until either the first mismatch is found or an entire match has been verified. This comparison can happen right-to-left, left-to-right, or in an arbitrary order that may depend on $p$. In our analysis, we focus on the right-to-left case for concreteness, but the modifications for the other cases are straightforward. Therefore, the cost of window $w$ is

$$\xi^p(w) = \begin{cases} m & \text{if } p = w, \\ \min\{i : 1 \leq i \leq m,\ p[m-i] \neq w[m-i]\} & \text{otherwise.} \end{cases}$$

In any case, the rightmost window character $a$ is used to determine how far the window can be shifted for the next iteration. The shift-function ensures that no match can be missed by moving the window such that $a$ becomes aligned to the rightmost $a$ in $p$ (not considering the last position). If $a$ does not occur in $p$ (or only at the last position), it is safe to shift by $m$ positions. Formally, we define

$$\begin{aligned} right^p(a) &:= \max\left[\{i \in \{0,\ldots,m-2\} : p[i] = a\} \cup \{-1\}\right], \\ \texttt{shift}[a] &:= (m-1) - right^p(a),\ \text{assuming } p \text{ fixed,} \\ shift^p(w) &:= \texttt{shift}[w[m-1]]. \end{aligned}$$

For concreteness, we state Horspool's algorithm and how we count text character accesses as pseudocode in Algorithm 1. Note that after a shift, even when we know that $a$ now matches its corresponding pattern character, the corresponding position is compared again and counts as a text access. Otherwise the additional bookkeeping would make the algorithm more complicated; this is not worth the effort in practice. The lookup in the `shift`-table does not count as an additional access, since we can remember $\texttt{shift}[a]$ as soon as the last window character has been read.

The main advantage of the Horspool algorithm is its simplicity. Especially, a window's shift value depends only on its last character, and its cost is easily computed from the number of consecutive matching characters at its right end. The Horspool algorithm does not require any advanced data structure and can be implemented in a few lines of code.

## 2.2 Backward (Nondeterministic) DAWG Matching, B(N)DM

The main idea of the BDM algorithm is to build a deterministic finite automaton (in this case, a suffix automaton, which is a directed acyclic word graph or DAWG) that recognizes all substrings of the reversed pattern, accepts all suffixes of the reversed pattern (including the empty suffix), and enters a FAIL state if a string has been read that is not a substring of the reversed pattern.

The suffix automaton processes the window right-to-left. As long as the FAIL state has not been reached, we have read a substring of the reversed pattern. If we are in an accepting state, we have even found a suffix of the reversed pattern (i.e., a prefix of $p$).

**Algorithm 1** Horspool-with-Cost

**Input:** text $s \in \Sigma^*$, pattern $p \in \Sigma^m$
**Output:** pair (number $occ$ of occurrences of $p$ in $s$, number $cost$ of accesses to $s$)
1: pre-compute table $\mathtt{shift}[a]$ for all $a \in \Sigma$
2: $(occ, cost) \leftarrow (0, 0)$
3: $t \leftarrow m - 1$
4: **while** $t < |s|$ **do**
5:    $i \leftarrow 0$
6:    **while** $i < m$ **do**
7:       $cost \leftarrow cost + 1$
8:       **if** $s[t - i] \neq p[(m - 1) - i]$ **then break**
9:       $i \leftarrow i + 1$
10:    **if** $i = m$ **then** $occ \leftarrow occ + 1$
11:    $t \leftarrow t + \mathtt{shift}[s[t]]$
12: **return** $(occ, cost)$

Whenever this happens before we have read $m$ characters, the last such event marks the next potential window start that may contain a match with $p$, and hence determines the shift. When we are in an accepting state after reading $m$ characters, we have found a match, but this does not influence the shift.

So, $\xi^p(w)$ is the number of characters read when entering FAIL, or $m$ if $p = w$. Let $I^p(w) \subseteq \{0, \ldots, m-1\}$ be the set defined by $i \in I^p(w)$ if and only if the suffix automaton of $\overleftarrow{p}$ is in an accepting state after reading $i$ characters of $w$. Then

$$shift^p(w) = \min \left\{ m - i \,\middle|\, i \in I^p(w) \right\}.$$

Note that $I^p(w)$ is never empty as the suffix automaton accepts the empty string and, thus, $0 \in I^p(w)$ for all windows $w$.

The advantage of BDM are long shifts, but its main disadvantage is the necessary construction of the suffix automaton, which is possible in $O(m)$ time via the suffix tree of the reversed pattern, but too expensive in practice to compete with other algorithms unless the search text is extremely long.

Constructing a nondeterministic finite automaton (NFA) instead of the deterministic suffix automaton is much simpler. However, processing a text character then does not take constant, but $O(m)$ time. However, the NFA can be efficiently simulated with bit-parallel operations such that processing a text character takes $O(m/W)$ time, where $W$ is the machine word size. For many patterns in practice, this is as good as $O(1)$. The resulting algorithm is then called BNDM.

From the "text character accesses" analysis point of view, BDM and BNDM are equivalent, as they have the same shift and cost functions.

## 2.3 Backward Oracle Matching, BOM

BOM is similar to BDM, but the suffix automaton of the reversed pattern is replaced by a simpler deterministic automaton, the factor oracle (rarely also called suffix oracle). It may recognize (accept) more strings than substrings (suffixes) of the reversed pattern, but is much easier to construct. It still guarantees that, once the FAIL state is reached, the sequence of read characters is not a substring of the reversed pattern.

The cost and shift functions are defined as for BDM, but based on the oracle. We refer to [16] for the construction details and further properties of the oracle. By construction, BOM never gives longer shifts than B(N)DM. The main advantage of BOM over BDM is reduced space usage and preprocessing time; the factor oracle only has $m + 1$ states and can be constructed faster than a suffix automaton.

# 3 Deterministic Arithmetic Automata

In this section, we introduce deterministic arithmetic automata (DAAs). They extend ordinary deterministic finite automata (DFAs) by performing a computation while one moves from state to state. Even though DAAs can be shown to be formally equivalent to families of DFAs on an appropriately defined larger state space, they are a useful concept before introducing probabilistic arithmetic automata (PAAs) and allow us to construct PAAs for the analysis of pattern matching algorithms in a simpler way.

**Definition 1** (Deterministic Arithmetic Automaton, DAA)**.** A *deterministic arithmetic automaton* is a tuple

$$\mathcal{D} = \big(\mathcal{Q}, q_0, \Sigma, \delta, \mathcal{V}, v_0, \mathcal{E}, (\eta_q)_{q \in \mathcal{Q}}, (\theta_q)_{q \in \mathcal{Q}}\big),$$

where $\mathcal{Q}$ is a finite set of states, $q_0 \in \mathcal{Q}$ is the start state, $\Sigma$ is a finite alphabet, $\delta : \mathcal{Q} \times \Sigma \to \mathcal{Q}$ is a transition function, $\mathcal{V}$ is a finite or countable set of values, $v_0 \in \mathcal{V}$ is called the start value, $\mathcal{E}$ is a finite set of emissions, $\eta_q \in \mathcal{E}$ is the emission associated to state $q$, and $\theta_q : \mathcal{V} \times \mathcal{E} \to \mathcal{V}$ is a binary operation associated to state $q$.

Informally, a DAA starts with the state-value pair $(q_0, v_0)$ and reads a sequence of symbols from $\Sigma$. Being in state $q$ with value $v$, upon reading $\sigma \in \Sigma$, the DAA performs a state transition to $q' := \delta(q, \sigma)$ and updates the value to $v' := \theta_{q'}(v, \eta_{q'})$ using the operation and emission of the new state $q'$.

Further, we define the associated joint transition function

$$\hat{\delta} : (\mathcal{Q} \times \mathcal{V}) \times \Sigma \to (\mathcal{Q} \times \mathcal{V}), \qquad \hat{\delta}\big((q, v), \sigma\big) := \big(\delta(q, \sigma), \theta_{\delta(q, \sigma)}(v, \eta_{\delta(q, \sigma)})\big).$$

As usual, we extend the definition of $\hat{\delta}$ inductively from $\Sigma$ to $\Sigma^*$ in its second argument by $\hat{\delta}\big((q, v), \varepsilon\big) := (q, v)$ for the empty string $\varepsilon$ and $\hat{\delta}\big((q, v), x\sigma\big) := \hat{\delta}\big(\hat{\delta}((q, v), x), \sigma\big)$ for all $x \in \Sigma^*$ and $\sigma \in \Sigma$. When $\hat{\delta}\big((q_0, v_0), s\big) = (q, v)$ for some $q \in \mathcal{Q}$ and $s \in \Sigma^*$, we say that $\mathcal{D}$ computes value $v$ for input $s$ and define $value_{\mathcal{D}}(s) := v$.

For each state $q$, the emission $\eta_q$ is fixed and could be dropped from the definition of DAAs. In fact, one could also dispense with values and operations entirely and define a DFA over state space $\mathcal{Q} \times \mathcal{V}$, performing the same operations as a DAA. However, we intentionally include values, operations, and emissions to emphasize the connection to PAAs (which are defined in Section 5).

As a simple example for a DAA, take a standard DFA $(\mathcal{Q}, q_0, \Sigma, \delta, F)$ with $F \subset \mathcal{Q}$ being a set of final (or accepting) states. To obtain a DAA that counts how many times the DFA visits an accepting state when reading $x \in \Sigma^*$, let $\mathcal{E} := \{0, 1\}$ and define $\eta_q := 1$ if $q \in F$, and $\eta_q := 0$ otherwise. Further define $\mathcal{V} = \mathbb{N}$ with $v_0 := 0$, and let the operation in each state be the usual addition: $\theta_q(v, e) := v + e$ for all $q$. Then $value_{\mathcal{D}}(x)$ is the desired count.

# 4 Constructing DAAs for Pattern Matching Analysis

For a given algorithm and pattern $p \in \Sigma^m$ with known shift and cost functions, $shift^p : \Sigma^m \to \{1, \ldots, m\}$, $w \mapsto shift^p(w)$ and $\xi^p : \Sigma^m \to \mathbb{N}$, $w \mapsto \xi^p(w)$, we construct a DAA that upon reading a text $s \in \Sigma^*$ computes the total cost, defined as the sum of costs over all examined windows. (Which windows are examined depends of course on the shift values of previously examined windows.) Slightly abusing notation, we write $\xi^p(s)$ for the total cost incurred on $s$.

While different constructions are possible (see also [15]), the construction presented here has the advantage that it is simple to describe and implement and processes only one text character at a time. This property allows the construction of a product DAA that directly compares two algorithms as detailed in Section 6.

We define a DAA by

- $\mathcal{Q} := \Sigma^m \times \{0, \ldots, m\}$,

- $q_0 := (p, m)$,

- $\mathcal{V} := \mathbb{N}$,

- $v_0 := 0$,

- $\mathcal{E} := \{1, \ldots, m\}$,

- $\eta_{(w,x)} := \begin{cases} 0 & \text{if } x > 0, \\ \xi^p(w) & \text{if } x = 0, \end{cases}$

- $\theta_q : (v, e) \mapsto v + e$ for all $q \in \mathcal{Q}$ (addition),

- $\delta : \big((w, x), \sigma\big) \mapsto \begin{cases} (w'\sigma, \ x - 1) & \text{if } x > 0, \\ (w'\sigma, \ shift^p(w) - 1) & \text{if } x = 0, \end{cases}$
  where $w'$ is the length-$(m-1)$ suffix of $w$, i.e., $w' := w[1] \ldots w[m-1]$.

Informally, the state $q = (w, x)$ means that the last $m$ read characters spell $w$ and that $x$ more characters need to be read to get to the end of the current window. For the start state $(p, m)$, the component $p$ is arbitrary, as we need to read $m$ characters to reach the end of the first window. The value accumulates the cost of examined windows. Therefore, the operation is a simple addition in each state, and the emission of state $(w, x)$ specifies the cost to add. Consequently, the emission is zero if the state does not correspond to an examined window ($x > 0$), and the emission equals the window cost $\xi^p(w)$ if $x = 0$. The transition function $\delta$ specifies how to move from one state to the next when reading the next text character $\sigma \in \Sigma$: In any case, the window content is updated by forgetting the first character and appending the read $\sigma$. If the end of the current window has not been reached ($x > 0$), the counter $x$ is decremented. Otherwise, the window's shift value is used to compute the number of characters till the next window aligns.

**Theorem 1.** With the DAA $\mathcal{D}$ constructed as above, $value_{\mathcal{D}}(s) = \xi^p(s)$ for all $s \in \Sigma^*$.

*Proof.* The total cost $\xi^p(s)$ can be written as the sum of costs of all processed windows: $\xi^p(s) = \sum_{i \in I_s} \xi^p(s[i - m + 1 \ldots i])$, where $I_s$ is the set of indices giving the processed windows, i.e. $I_s \subset \{m - 1, \ldots, |s| - 1\}$ such that

$$i \in I_s \quad :\Longleftrightarrow \quad i = m - 1 \quad \text{or} \quad \exists j \in I_s : i = j + shift^p(s[j - m + 1 \ldots j]).$$

We have to prove that the DAA computes this value for $s \in \Sigma^*$.

Let $(w_i, x_i)$ be the DAA state active after reading $s[..i]$. Observe that the transition function $\delta$ ensures that the $w_i$-component of $(w_i, x_i)$ reflects the rightmost length-$m$ window of $s[..i]$, which can immediately be verified inductively. Thus, the emission on reading the last character $s[i]$ of $s[..i]$ with $i \geq m - 1$ is, by definition of $\eta_{(w_i, x_i)}$, either $\xi^p(s[i - m + 1 \ldots i])$ or zero, depending on the second component of $(w_i, x_i)$. As the operation is an addition for all states, $value_{\mathcal{D}}(s) = \sum_{i \in I'_s} \xi^p(s[i - m + 1 \ldots i])$ for

$$I'_s := \big\{ i \in \{0, \ldots, |s| - 1\} : x_i = 0 \big\}.$$

It remains to show that $I_s = I'_s$. To this end, note that by $\delta$, we have $x_{i+1} = x_i - 1$ if $x_{i+1} > 0$ and $x_{i+1} = shift^p(w_i) - 1$ if $x_{i+1} = 0$. As $q_0 = (p, m)$, it follows that $m - 1 \in I'_s$. Using $w_i = s[i - m + 1 \ldots i]$ for $i \geq m - 1$, we conclude that whenever $x_i = 0$, it follows that $x_j = 0$ for $j = i + shift^p(s[i - m + 1 \ldots i])$ and that $x_{j'} > 0$ for $i < j' < j$. Hence we obtain that $i \in I'_s$ implies that $i + shift^p(s[i - m + 1 \ldots i]) \in I'_s$ and $i + k \notin I'_s$ for $0 < k < shift^p(s[i - m + 1 \ldots i])$, which completes the proof. $\square$

**DAA Minimization** The size of the constructed DAA's state space depends exponentially on the pattern length, making the application for long patterns infeasible in practice. However, depending on the particular circumstances (i.e., algorithm and pattern analyzed), the constructed DAA can often be substantially reduced by applying a modified version of Hopcroft's algorithm for DFA minimization [7].

Hopcroft's algorithm minimizes a DFA in $O(|\mathcal{Q}| \log |\mathcal{Q}|)$ time by iteratively refining a partition of the state set. In the beginning, all states are partitioned into two distinct

sets: one containing the accepting stats, and the other containing the non-accepting states. This partition is iteratively refined whenever a reason for non-equivalence of two states in the same set is found. Upon termination, the states are partitioned into sets of equivalent states. Refer to [10] for an in-depth explanation of Hopcroft's algorithm.

The algorithm can straightforwardly be adapted to minimize DAAs by choosing the initial state set partition appropriately. In our case, each DAA state is associated with the same operation. The only differences in state's behavior thus stem from different emissions. Therefore, Hopcroft's algorithm can be initialized by the partition induced by the emissions and then continued as usual.

As we exemplify in Section 7, this leads to a considerable reduction of the number of states.

# 5 Probabilistic Arithmetic Automata

This section introduces finite-memory random text models and explains how to construct a *probabilistic arithmetic automaton* (PAA) from a (minimized) DAA and a random text model. PAAs were introduced in [13], where they are used to compute pattern occurrence count distributions. Other applications in biological sequence analysis include the exact computation of p-values of sequence motifs [14], and the determination of seed sensitivity for pairwise sequence alignment algorithms based on filtering [6].

## 5.1 Random Text Models

Given an alphabet $\Sigma$, a random text is a stochastic process $(S_t)_{t \in \mathbb{N}_0}$, where each $S_t$ takes values in $\Sigma$. A text model $\mathbb{P}$ is a probability measure assigning probabilities to (sets of) strings. It is given by (consistently) specifying the probabilities $\mathbb{P}(S_0 \ldots S_{|s|-1} = s)$ for all $s \in \Sigma^*$. We only consider finite-memory models in this article which are formalized in the following definition.

**Definition 2** (Finite-memory text model)**.** A finite-memory text model is a tuple $(\mathcal{C}, c_0, \Sigma, \varphi)$, where $\mathcal{C}$ is a finite state space (called *context space*), $c_0 \in \mathcal{C}$ a start context, $\Sigma$ an alphabet, and $\varphi : \mathcal{C} \times \Sigma \times \mathcal{C} \to [0, 1]$ a transition function with $\sum_{\sigma \in \Sigma, c' \in \mathcal{C}} \varphi(c, \sigma, c') = 1$ for all $c \in \mathcal{C}$. The random variable giving the text model state after $t$ steps is denoted $C_t$ with $C_0 :\equiv c_0$. A probability measure is now induced by stipulating

$$\mathbb{P}(S_0 \ldots S_{n-1} = s, C_1 = c_1, \ldots, C_n = c_n) := \prod_{i=0}^{n-1} \varphi(c_i, s[i], c_{i+1})$$

for all $n \in \mathbb{N}_0$, $s \in \Sigma^n$, and $(c_1, \ldots, c_n) \in \mathcal{C}^n$.

The idea is that the model given by $(\mathcal{C}, c_0, \Sigma, \varphi)$ generates a random text by moving from context to context and emitting a character at each transition, where $\varphi(c, \sigma, c')$ is the probability of moving from context $c$ to context $c'$ and thereby generating the letter $\sigma$.

9

Note that the probability $\mathbb{P}(S_0 \ldots S_{|s|-1} = s)$ is obtained by marginalization over all context sequences that generate $s$. This can be efficiently done, using the decomposition of the following lemma.

**Lemma 1.** Let $(\mathcal{C}, c_0, \Sigma, \varphi)$ be a finite-memory text model. Then,

$$\mathbb{P}(S_0 \ldots S_n = s\sigma, C_{n+1} = c) = \sum_{c' \in \mathcal{C}} \mathbb{P}(S_0 \ldots S_{n-1} = s, C_n = c') \cdot \varphi(c', \sigma, c)$$

for all $n \in \mathbb{N}_0$, $s \in \Sigma^n$, $\sigma \in \Sigma$ and $c \in \mathcal{C}$.

*Proof.* We have

$$\begin{aligned}
&\mathbb{P}(S_0 \ldots S_n = s\sigma, C_{n+1} = c) \\
&= \sum_{c_1, \ldots, c_n} \mathbb{P}(S_0 \ldots S_n = s\sigma, C_1 = c_1, \ldots, C_n = c_n, C_{n+1} = c) \\
&= \sum_{c_1, \ldots, c_n} \prod_{i=0}^{n-1} \varphi(c_i, s[i], c_{i+1}) \cdot \varphi(c_n, \sigma, c) \\
&= \sum_{c_n \in \mathcal{C}} \left( \sum_{c_1, \ldots, c_{n-1}} \prod_{i=0}^{n-1} \varphi(c_i, s[i], c_{i+1}) \right) \cdot \varphi(c_n, \sigma, c) \\
&= \sum_{c_n \in \mathcal{C}} \mathbb{P}(S_0 \ldots S_{n-1} = s, C_n = c_n) \cdot \varphi(c_n, \sigma, c) \,.
\end{aligned}$$

Renaming $c_n$ to $c'$ yields the claimed result. $\qquad \square$

Similar text models are used in [11], where they a called probability transducers. In the following, we refer to a finite-memory text model $(\mathcal{C}, c_0, \Sigma, \varphi)$ simply as text model, as all text models considered in this article are special cases of Definition 2.

For an i.i.d. model, we set $\mathcal{C} = \{\varepsilon\}$ and $\varphi(\varepsilon, \sigma, \varepsilon) = p_\sigma$ for each $\sigma \in \Sigma$, where $p_\sigma$ is the occurrence probability of letter $\sigma$ (and $\varepsilon$ may be interpreted as an empty context). For a Markovian text model of order $r$, the distribution of the next character depends on the $r$ preceding characters (fewer at the beginning); thus $\mathcal{C} := \bigcup_{i=0}^{r} \Sigma^i$. This notion of text models also covers variable order Markov chains as introduced in [17], which can be converted into equivalent models of fixed order. Text models as defined above have the same expressive power as character-emitting HMMs, that means, they allow to construct the same probability distributions.

## 5.2 Basic PAA Concepts

Probabilistic arithmetic automata (PAAs), as introduced in [13], are a generic concept useful to model probabilistic chains of operations. In this section, we sum up the definition and basic recurrences needed in this article.

**Definition 3** (Probabilistic Arithmetic Automaton, PAA). A *probabilistic arithmetic automaton* is a tuple $\mathcal{P} = \left( \mathcal{Q}, q_0, T, \mathcal{V}, v_0, \mathcal{E}, \mu = (\mu_q)_{q \in \mathcal{Q}}, \theta = (\theta_q)_{q \in \mathcal{Q}} \right)$, where $\mathcal{Q}$, $q_0$, $\mathcal{V}$, $v_0$, $\mathcal{E}$ and $\theta$ have the same meaning as for a DAA, each $\mu_q$ is a state-specific probability distribution on the emissions $\mathcal{E}$, and $T : \mathcal{Q} \times \mathcal{Q} \to [0, 1]$ is a transition function, such that $T(q, q')$ gives the probability of a transition from state $q$ to state $q'$, i.e. $\left( T(q, q') \right)_{q, q' \in \mathcal{Q}}$ is a stochastic matrix.

A PAA induces three stochastic processes: (1) the state process $(Q_t)_{t \in \mathbb{N}}$ with values in $\mathcal{Q}$, (2) the emission process $(E_t)_{t \in \mathbb{N}}$ with values in $\mathcal{E}$, and (3) the value process $(V_t)_{t \in \mathbb{N}}$ with values in $\mathcal{V}$ such that $\mathcal{V}_0 :\equiv v_0$ and $\mathcal{V}_t := \theta_{Q_t} (V_{t-1}, E_t)$.

We now restate the PAA recurrences from [13] to compute the state-value distribution after $t$ steps. For the sake of a shorter notation, we define $f_t(q, v) := \mathbb{P}(Q_t = q, V_t = v)$. Since we are generally only interested in the value distribution, note that it can be obtained by marginalization: $\mathbb{P}(V_t = v) = \sum_{q \in \mathcal{Q}} f_t(q, v)$.

**Lemma 2** (State-value recurrence, [13]).
The state-value distribution is given by $f_0(q, v) = 1$ if $q = q_0$ and $v = v_0$, and $f_0(q, v) = 0$ otherwise. For $t \geq 0$,

$$f_{t+1}(q, v) = \sum_{q' \in \mathcal{Q}} \sum_{(v', e) \in \theta_q^{-1}(v)} f_t(q', v') \cdot T(q', q) \cdot \mu_q(e), \tag{1}$$

where $\theta_q^{-1}(v)$ denotes the inverse image set of $v$ under $\theta_q$.

The recurrence in Lemma 2 resembles the Forward recurrences known from HMMs.

Note that the range of $V_t$ is finite for each $t$, even when $\mathcal{V}$ is infinite, as $V_t$ is a function of the states and emissions up to time $t$, and state set $\mathcal{Q}$ and emission set $\mathcal{E}$ are finite. We define $\mathcal{V}_t := \text{range } V_t$ and $\vartheta_n := \max_{0 \leq t \leq n} |\mathcal{V}_t|$. Clearly $\vartheta_n \leq (|\mathcal{Q}| \cdot |\mathcal{E}|)^n$. Therefore all actual computations are on finite sets. When analyzing the number of character accesses of a pattern matching algorithm, we have $\mathcal{V}_t \subset \{0, \ldots, m(n - m + 1)\}$, as at most $(n - m + 1)$ search windows are processed, each causing at most $m$ character accesses. Thus, $\vartheta_n \in O(n \cdot m)$.

## 5.3 Constructing a PAA from a DAA and a Text Model

We now formally state how to combine a DAA and a text model into a PAA that allows us to compute the distribution of values produced by the DAA when processing a random text.

**Lemma 3** (DAA + Text model $\to$ PAA). Let $(\mathcal{C}, c_0, \Sigma, \varphi)$ be a text model and $\mathcal{D} = \left( \mathcal{Q}^{\mathcal{D}}, q_0^{\mathcal{D}}, \Sigma, \delta, \mathcal{V}, v_0, \mathcal{E}, (\eta_q)_{q \in \mathcal{Q}^{\mathcal{D}}}, (\theta_q^{\mathcal{D}})_{q \in \mathcal{Q}^{\mathcal{D}}} \right)$ be a DAA. Then, define

- a state space $\mathcal{Q} := \mathcal{Q}^{\mathcal{D}} \times \mathcal{C}$,

- a start state $q_0 := (q_0^{\mathcal{D}}, c_0)$,

- transition probabilities

$$T\big((q,c),(q',c')\big) := \sum_{\sigma \in \Sigma:\, \delta(q,\sigma)=q'} \varphi(c,\sigma,c'), \tag{2}$$

- (deterministic) emission probability vectors

$$\mu_{(q,c)}(e) := \begin{cases} 1 & \text{if } e = \eta_q\,, \\ 0 & \text{otherwise}\,, \end{cases}$$

for all $(q,c) \in \mathcal{Q}$.

- operations $\theta_{(q,c)}(v,e) := \theta_q^{\mathcal{D}}(v,e)$ for all $(q,c) \in \mathcal{Q}$.

Then, $\mathcal{P} = \big(\mathcal{Q}, q_0, T, \mathcal{V}, v_0, \mathcal{E}, \mu = (\mu_q)_{q \in \mathcal{Q}}, \theta = (\theta_q)_{q \in \mathcal{Q}}\big)$ is a PAA with

$$\mathcal{L}(V_t) = \mathcal{L}\big(value_{\mathcal{D}}(S_0 \ldots S_{t-1})\big),$$

for all $t \in \mathbb{N}_0$, where $S$ is a random text according to the text model $(\mathcal{C}, c_0, \Sigma, \varphi)$. States having zero probability of being reached from $q_0$ may be omitted from $\mathcal{Q}$ and $T$. For such a PAA, the value distribution $\mathcal{L}(V_n)$ can be computed with $O(n \cdot |\mathcal{Q}^{\mathcal{D}}| \cdot |\mathcal{C}|^2 \cdot |\Sigma| \cdot \vartheta_n)$ operations using $O(|\mathcal{Q}^{\mathcal{D}}| \cdot |\mathcal{C}| \cdot \vartheta_n)$ space. If for all $c \in \mathcal{C}$ and $\sigma \in \Sigma$, there exists at most one $c' \in \mathcal{C}$ such that $\varphi(c,\sigma,c') > 0$, then the runtime is bounded by $O(n \cdot |\mathcal{Q}^{\mathcal{D}}| \cdot |\mathcal{C}| \cdot |\Sigma| \cdot \vartheta_n)$.

*Proof.* $\mathcal{P}$ is a PAA by Definition 3. As in Section 5.2, we define $f_t(q,v) := \mathbb{P}(Q_t = q, V_t = v)$. Iverson brackets are written $[\![\cdot]\!]$, i.e. $[\![A]\!] = 1$ if the statement $A$ is true and $[\![A]\!] = 0$ otherwise.

To prove $\mathcal{L}(V_t) = \mathcal{L}\big(value_{\mathcal{D}}(S_0 \ldots S_{t-1})\big)$, we show that

$$f_t\big((q^{\mathcal{D}},c),v\big) = \sum_{s \in \Sigma^t} [\![\hat{\delta}\big((q_0^{\mathcal{D}},v_0),s\big) = (q^{\mathcal{D}},v)]\!] \cdot \mathbb{P}(S_0 \ldots S_{t-1} = s, C_t = c) \tag{3}$$

for all $q^{\mathcal{D}} \in \mathcal{Q}^{\mathcal{D}}$, $c \in \mathcal{C}$, $v \in \mathcal{V}$, and $t \in \mathbb{N}_0$. For $t = 0$, Equation (3) is correct by definitions of PAAs, DAAs and text models. For $t > 0$ we prove it inductively.

Assume (3) to be correct for all $t'$ with $0 \le t' < t$. Then

$$f_t\big(\underbrace{(q^{\mathcal{D}}, c)}_{=:q}, v\big) \tag{4}$$

$$= \sum_{q' \in \mathcal{Q}} \sum_{(v', e) \in \theta_q^{-1}(v)} f_{t-1}(q', v') \cdot T(q', q) \cdot \mu_q(e) \tag{5}$$

$$= \sum_{q' \in \mathcal{Q}} \sum_{(v', e) \in \mathcal{V} \times \mathcal{E}} [\![\theta_{q^{\mathcal{D}}}^{\mathcal{D}}(v', e) = v]\!] \cdot f_{t-1}(q', v') \cdot T(q', q) \cdot [\![\eta_{q^{\mathcal{D}}} = e]\!] \tag{6}$$

$$= \sum_{q'^{\mathcal{D}} \in \mathcal{Q}^{\mathcal{D}}} \sum_{c' \in \mathcal{C}} \sum_{(v', e) \in \mathcal{V} \times \mathcal{E}} [\![\theta_{q^{\mathcal{D}}}^{\mathcal{D}}(v', e) = v]\!] \cdot [\![\eta_{q^{\mathcal{D}}} = e]\!] \cdot f_{t-1}(q', v')$$
$$\cdot \sum_{\sigma \in \Sigma} [\![\delta(q'^{\mathcal{D}}, \sigma) = q^{\mathcal{D}}]\!] \cdot \varphi(c', \sigma, c) \tag{7}$$

$$= \sum_{s \in \Sigma^{t-1}} \sum_{\sigma \in \Sigma} \sum_{q'^{\mathcal{D}} \in \mathcal{Q}^{\mathcal{D}}} \sum_{c' \in \mathcal{C}} \sum_{(v', e) \in \mathcal{V} \times \mathcal{E}} [\![\theta_{q^{\mathcal{D}}}^{\mathcal{D}}(v', e) = v]\!] \cdot [\![\eta_{q^{\mathcal{D}}} = e]\!]$$
$$\cdot [\![\delta(q'^{\mathcal{D}}, \sigma) = q^{\mathcal{D}}]\!] \cdot [\![\hat{\delta}\big((q_0^{\mathcal{D}}, v_0), s\big) = (q'^{\mathcal{D}}, v')]\!]$$
$$\cdot \mathbb{P}(S_0 \dots S_{t-2} = s, C^{t-1} = c') \cdot \varphi(c', \sigma, c) \tag{8}$$

$$= \sum_{s\sigma \in \Sigma^t} \sum_{q'^{\mathcal{D}} \in \mathcal{Q}^{\mathcal{D}}} \sum_{(v', e) \in \mathcal{V} \times \mathcal{E}} [\![\theta_{q^{\mathcal{D}}}^{\mathcal{D}}(v', e) = v]\!] \cdot [\![\eta_{q^{\mathcal{D}}} = e]\!] \cdot [\![\hat{\delta}\big((q_0^{\mathcal{D}}, v_0), s\big) = (q'^{\mathcal{D}}, v')]\!]$$
$$\cdot [\![\delta(q'^{\mathcal{D}}, \sigma) = q^{\mathcal{D}}]\!] \cdot \mathbb{P}(S_0 \dots S_{t-1} = s\sigma, C_t = c) \tag{9}$$

$$= \sum_{s\sigma \in \Sigma^t} [\![\hat{\delta}\big((q_0^{\mathcal{D}}, v_0), s\sigma\big) = (q^{\mathcal{D}}, v)]\!] \cdot \mathbb{P}(S_0 \dots S_{t-1} = s\sigma, C_t = c) \tag{10}$$

In the above derivation, step (4)→(5) follows from (1). Step (5)→(6) follows from the definitions of $\theta_q$ and $\mu_q$. Step (6)→(7) uses the definitions of $T$ and $\mathcal{Q}$ in Lemma 3. Step (7)→(8) uses the induction assumption. Step (8)→(9) uses Lemma 1. The final step (9)→(10) follows by combining the four Iverson brackets summed over $q'^{\mathcal{D}}$ and $(v', e)$ into a single Iverson bracket.

To compute the table $f_n$ containing $f_n(q, v)$ for all $q \in \mathcal{Q}$ and $v \in \mathcal{V}$, we start with $f_0$ and perform $n$ update steps. The given runtime bounds can be verified by considering a "push" algorithm: When computing $f_{t+1}$, we initialize the table with zeros and iterate over all $q \in \mathcal{Q}$, $v \in \mathcal{V}$ and $q' \in \{q'' \in \mathcal{Q} : T(q, q'') > 0\}$; for each combination of $q$, $v$, and $q'$, we add $f_t(q, v) \cdot T(q, q')$ to $f_{t+1}\big(q', \theta_{q'}(v, \eta_{q'})\big)$. $\square$

As a direct consequence of the above lemma and of the DAA construction from Section 4, we arrive at our main theorem.

**Theorem 2.** Given a finite-memory text model $(\mathcal{C}, c_0, \Sigma, \varphi)$, a window-based pattern matching algorithm $A$, a pattern $p$ with $|p| = m$, and the functions $shift^{A,p}$ and $\xi^{A,p}$, the cost distribution $\mathcal{L}(X_n^{A,p})$ can be computed using $O(n^2 \cdot m \cdot |\mathcal{Q}^{\mathcal{D}}| \cdot |\mathcal{C}|^2 \cdot |\Sigma|)$ time and $O(|\mathcal{Q}^{\mathcal{D}}| \cdot |\mathcal{C}| \cdot n \cdot m)$ space. Since $|\mathcal{Q}^{\mathcal{D}}|$ is bounded by $O(m \cdot \Sigma^m)$, the computation uses $O(n^2 \cdot m^2 \cdot \Sigma^{m+1} \cdot |\mathcal{C}|^2)$ time and $O(m^2 \cdot \Sigma^m \cdot |\mathcal{C}| \cdot n)$ space. If for all $c \in \mathcal{C}$ and $\sigma \in \Sigma$,

there exists at most one $c' \in \mathcal{C}$ such that $\varphi(c, \sigma, c') > 0$, a factor of $|\mathcal{C}|$ can be dropped from the runtime bounds.

Applying DAA minimization before the PAA construction results in considerable speed-ups in practice. Alternatively, algorithm dependent construction schemes may be used to construct small automata. Tsai [20], for instance, gives algorithms to compute the asymptotic mean and variance of the number of comparisons used by Horspool's algorithm; for that, he constructs a Markov chain with $O(m^2)$ states. His construction can immediately be adapted to construct a DAA with $O(m^2)$ states.

# 6 Comparing Algorithms with Difference DAAs

Computing the cost distribution for two algorithms allows us to compare their performance characteristics. One natural question, however, cannot be answered by comparing these two (one-dimensional) distributions: What is the probability that algorithm $A$ needs more text accesses than algorithm $B$ to scan the same random text? The answer will depend on the correlation of algorithm performances: Do the same instances lead to long runtimes for both algorithms or are there instances that are easy for one algorithm but difficult for the other? This section answers these questions by constructing a PAA to compute the distribution of *cost differences* of two algorithms. That means, we calculate the probability that algorithm $A$ needs $v$ text accesses *more* than algorithm $B$ for all $v \in \mathbb{Z}$.

We start by giving a general construction of a DAA that computes the difference of the sum of emission of two given DAAs.

**Definition 4** (Difference DAA). Let $\mathcal{D}^1 = \left(\mathcal{Q}^1, q_0^1, \Sigma, \delta^1, \mathcal{V}^1, v_0^1, \mathcal{E}^1, (\eta_q^1)_{q \in \mathcal{Q}^1}, (\theta_q^1)_{q \in \mathcal{Q}^1}\right)$ and $\mathcal{D}^2 = \left(\mathcal{Q}^2, q_0^2, \Sigma, \delta^2, \mathcal{V}^2, v_0^2, \mathcal{E}^2, (\eta_q^2)_{q \in \mathcal{Q}^2}, (\theta_q^2)_{q \in \mathcal{Q}^2}\right)$ be DAAs given over the same alphabet $\Sigma$ with $\mathcal{V}^1 = \mathcal{V}^2 = \mathbb{N}$, $v_0^1 = v_0^2 = 0$, $\mathcal{E}^1, \mathcal{E}^2 \subset \mathbb{N}$, and all operations are additions of previous value and current emission. The *difference DAA* is defined as

$$DiffDAA(\mathcal{D}^1, \mathcal{D}^2) := (\mathcal{Q}, q_0, \Sigma, \delta, \mathcal{V}, v_0, \mathcal{E}, (\eta_q)_{q \in \mathcal{Q}}, (\theta_q)_{q \in \mathcal{Q}})$$

where

- $\mathcal{Q} := \mathcal{Q}^1 \times \mathcal{Q}^2$ and $q_0 := (q_0^1, q_0^2)$,

- $\mathcal{V} := \mathbb{Z}$ and $v_0 := 0$,

- $\mathcal{E} := \mathcal{E}^1 \times \mathcal{E}^2$ and $\eta_{(q^1, q^2)} := \left(\eta_{q^1}^1, \eta_{q^2}^2\right)$,

- $\delta : \left((q^1, q^2), \sigma\right) \mapsto \left(\delta^1(q^1, \sigma), \delta^2(q^2, \sigma)\right)$,

- $\theta_q : \left(v, (e^1, e^2)\right) \mapsto v + e^1 - e^2$

**Lemma 4.** Let $\mathcal{D}^1$ and $\mathcal{D}^2$ be DAAs meeting the criteria given in Definition 4 and $\mathcal{D} := DiffDAA(\mathcal{D}^1, \mathcal{D}^2)$. Then,

$$value_{\mathcal{D}}(s) = value_{\mathcal{D}^1}(s) - value_{\mathcal{D}^2}(s)$$

for all $s \in \Sigma^*$.

*Proof.* Follows directly from Definition 4. $\qquad\square$

Lemma 4 can now be applied to the DAAs constructed for the analysis of two algorithms as described in Section 4. Since the above construction builds the product of both state spaces, it is advisable to minimize both DAAs before generating the product. Furthermore, in an implementation, only reachable states of the product automaton need to be constructed. Before being used to build a PAA (by applying Lemma 3), the product DAA should again be minimized.

As discussed in Section 5.2, at most $m(n - m + 1)$ character accesses can result from scanning a text of length $n$ for a pattern of length $m$. Thus, the difference of costs for two algorithms lies in the range $\{-m(n - m + 1), \ldots, m(n - m + 1)\}$ and, hence, $\vartheta_n \in O(n \cdot m)$.

# 7 Case Studies

In Section 2, we considered three practically relevant algorithms, namely Horspool's algorithm, backward oracle matching (BOM), and backward (non)-deterministic DAWG matching (B(N)DM). Now, we compare the distributions of running time costs of these algorithms for several patterns. Figure 1 shows these distributions for the patterns `ATATAT` and `ACGTAC` for text lengths 100 and 500 under a uniform i.i.d. model on the DNA alphabet $\{A,C,G,T\}$. For text length 500, the distributions for Horspool and B(N)DM resemble the shape of normal distributions. In fact, for Horspool's algorithm it has been proven that the distribution is asymptotically normal [18]. For smaller text lengths (e.g. 100, as shown in left column of Figure 1), the distributions are less smooth than for longer texts. It is remarkable that for BOM we find zero probabilities with a fixed period. In all examples shown Figure 1 this period equals 7.

The probability that one pattern matching algorithm is faster than another depends on the pattern. Using the technique introduced in Section 6, we can quantify the strength of this effect. Figure 2 shows distributions of cost *differences* for different patterns and algorithms. That means, the probability that the first algorithm is faster is represented by the area under the curve left of zero. For the pattern `CGAAAA`, for example, there is a 55.6% probability that Horspool's algorithm needs fewer character accesses than B(N)DM in uniform i.i.d. texts of length 100, while for `ACGTAC`, the probability is only 0.18%.

Worth noting and perhaps surprising is the fact that there is a non-zero probability of BOM being faster than B(N)DM altough, $shift^{B(N)DM,p}(w) \geq shift^{BOM,p}(w)$ for all window contents $w$. The explanation, of course, is that a shorter (say, first) shift for
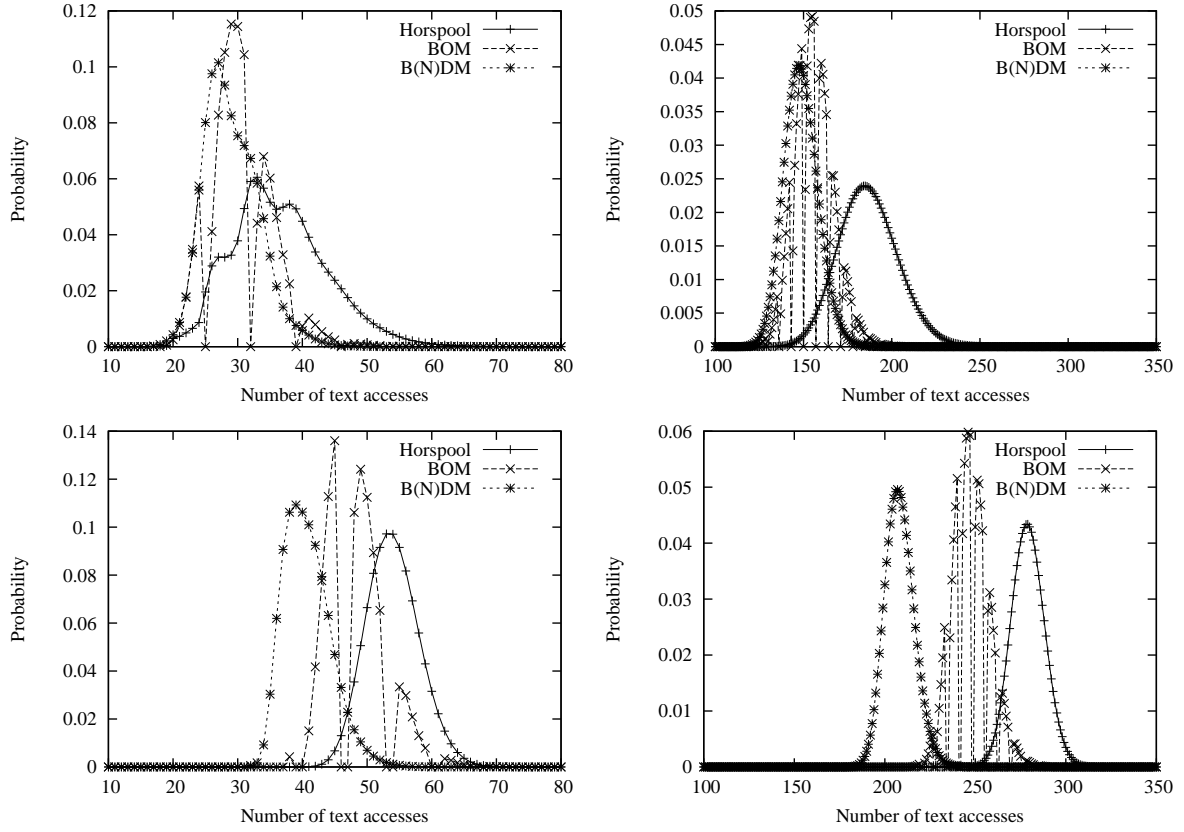
Figure 1: Exact distributions of character access counts for patterns `ATATAT` (top) and `ACGTAC` (bottom) for text length 100 (left) and text length 500 (right). An i.i.d. text model with uniform character distribution is used.

BOM leads to a different window content than for B(N)DM for the second window, which may have a larger shift value. This effect depends on the pattern: For the pattern `CAAAAA`, there is a 48.2% probability that BOM performs better than B(N)DM, while it is 6.2% for `ACGTAC`, again on texts of length 100.

To assess the effect of DAA minimization before constructing PAAs, we constructed minimized DAAs for all 21840 patterns of lengths 2 to 7 over $\Sigma = \{\texttt{A}, \texttt{C}, \texttt{G}, \texttt{T}\}$. The minimum, average, and maximum state counts are shown in Table 1. For length 6, Figure 3 contains a detailed histogram. These statistics show that construction and minimization as given in this article lead to smaller automata (and thus better runtimes) than the constructions given in the conference version of this article [15]. It may be conjectured that that minimal state spaces grows only polynomial with $m$ for all of these algorithms, as has been previously proven for the Horspool algorithm [20].

A JAVA implementation is available at `http://www.rahmannlab.de/software`. All algorithms were run on an Intel Core 2 Quad CPU at 2.66GHz. Computing the distributions shown in Figure 1 took 0.3 to 0.6 seconds for each distribution. Distributions of
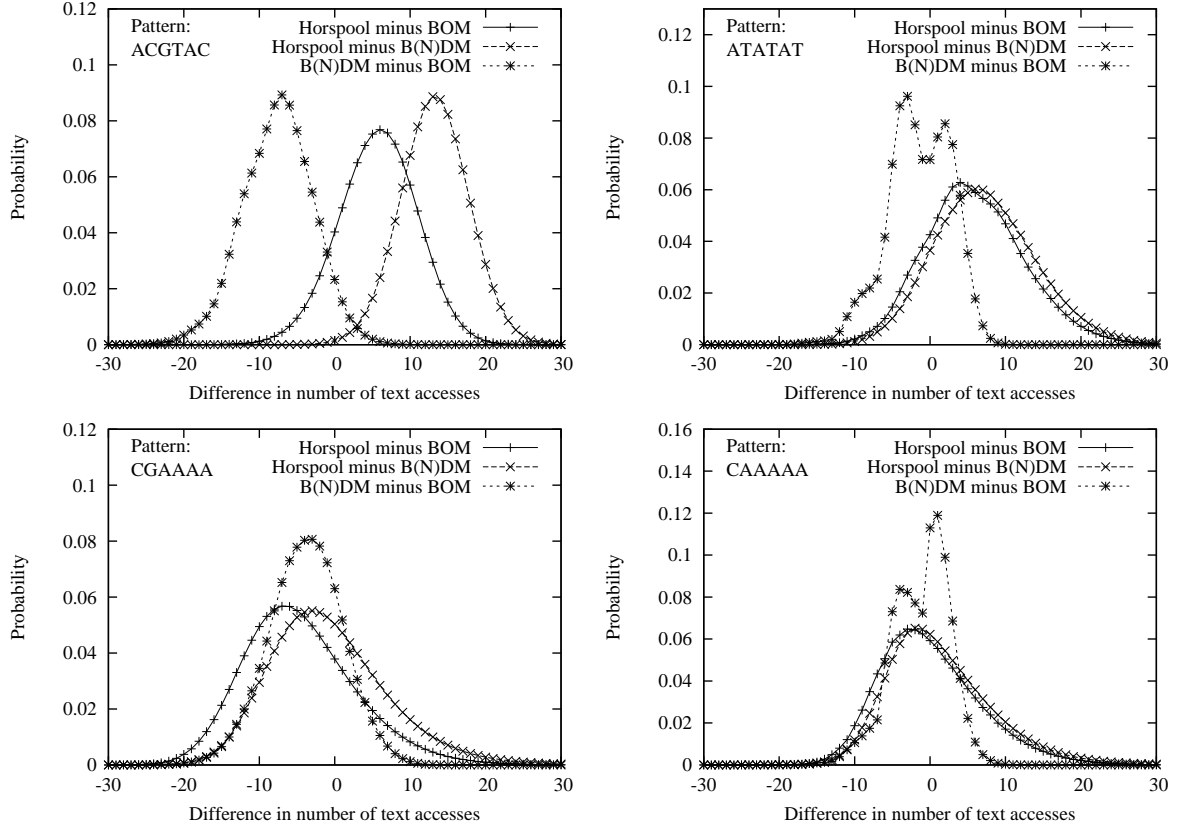
Figure 2: Exact distributions of differences in character access counts for different patterns using a uniform character distribution as text model and random texts of lengths 100.

differences as in Figure 2 were computed in 14 to 36 seconds.

# 8 Discussion

Using PAAs, we have shown how the exact distribution of the number of character accesses for window-based pattern matching algorithms can be computed. The framework is general enough to admit i.i.d. text models, Markovian text models of arbitrary order, and character-emitting hidden Markov models. The given construction results in an asymptotic runtime of $O(n^2 \cdot m \cdot |\mathcal{Q}^{\mathcal{D}}| \cdot |\mathcal{C}|^2 \cdot |\Sigma|)$. The number of DAA states $\mathcal{Q}^{\mathcal{D}}$ is $O(m \cdot \Sigma^m)$, but it can be considerably reduced by DAA minimization. The resulting PAA is smaller and therefore computing the cost distribution is much faster. If the pattern length $m$ is large, however, construction and minimization of the DAA itself pose a significant burden. It remains open if there exists an algorithm to construct the minimal automaton directly in general, i.e. using only $O(|\mathcal{Q}^{\mathcal{D}}_{\min}|)$ time. For Horspool's algorithm,

17

Table 1: Comparison of DAA sizes for all patterns of length $m$ over $\Sigma = \{A, C, G, T\}$.

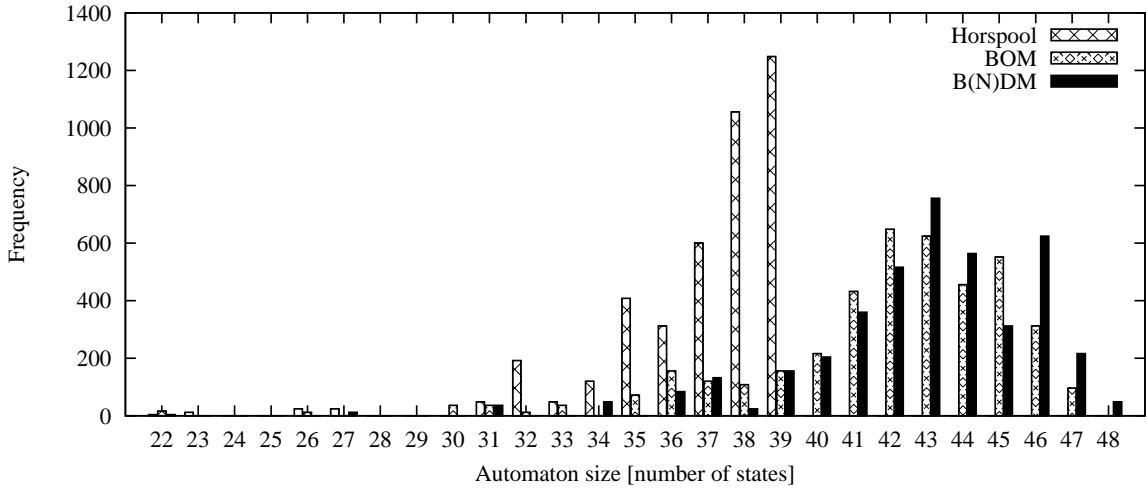| $m$ | States unminimized | States minimized (min./avg./max.) | | |
|:---:|:---:|:---:|:---:|:---:|
| | $|\Sigma|^m \cdot (m+1)$ | Horspool | BOM | B(N)DM |
| 2 | 48 | 4 / 4.8 / 5 | 4 / 4.0 / 4 | 4 / 4.8 / 5 |
| 3 | 256 | 7 / 8.3 / 9 | 7 / 8.3 / 9 | 7 / 9.6 / 10 |
| 4 | 1280 | 11 / 14.3 / 15 | 11 / 15.6 / 18 | 11 / 17.0 / 19 |
| 5 | 6144 | 16 / 23.6 / 25 | 16 / 26.5 / 30 | 16 / 27.9 / 31 |
| 6 | 28672 | 22 / 37.0 / 39 | 22 / 41.8 / 47 | 22 / 42.8 / 48 |
| 7 | 131072 | 29 / 55.2 / 58 | 29 / 62.4 / 70 | 29 / 62.6 / 70 |



Figure 3: Histogram on number of states of minimal DAAs over all patterns of length 6 over $\Sigma = \{A, C, G, T\}$.

an automaton with $O(m^2)$ states can be constructed by using the ideas from [20]. The construction, however, cannot be easily transferred to other algorithms and we are not aware of any similar results for BOM or B(N)DM.

In this work, we considered the most practically relevant algorithms. Exemplarily studying the cost distribution for some patterns showed that the algorithms performance is indeed non-negligibly influenced by the choice of the pattern. Especially the behavior of BOM deserves further attention: Its distribution of text character accesses features periodic zero probabilities, and unexpectedly, it may need fewer text accesses than B(N)DM on some patterns, although BOM's shift values are never better than B(N)DM's.

We focused on algorithms for single patterns, but the presented techniques also apply to algorithms to search for multiple patterns like the Wu-Manber algorithm [21] or set

backward oracle matching and multiple BNDM as given in [16]. A comparison of the resulting distributions could yield new insights into these algorithms as well.

Other metrics than text character accesses might be of interest and could be easily substituted; for example, just counting the number of windows by defining $\xi^p(w) = 1$ for all $w \in \Sigma^m$.

The given constructions allow us to analyze an algorithm's performance for each pattern individually. While this is desirable for detailed analysis, the cost distribution resulting from randomly choosing text *and* pattern would also be of interest.

# References

[1] C. Allauzen, M. Crochemore, and M. Raffinot. Efficient experimental string matching by weak factor recognition. In G. Goos, J. Hartmanis, and J. van Leeuwen, editors, *Proceedings of the 12th Annual Symposium on Combinatorial Pattern Matching (CPM)*, volume 2089 of *LNCS*, pages 51–72, 2001.

[2] R. A. Baeza-Yates, G. H. Gonnet, and M. Régnier. Analysis of Boyer-Moore-type string searching algorithms. In *SODA '90: Proceedings of the first annual ACM-SIAM symposium on Discrete algorithms*, pages 328–343. SIAM, 1990.

[3] R. A. Baeza-Yates and M. Régnier. Average running time of the Boyer-Moore-Horspool algorithm. *Theor. Comput. Sci.*, 92(1):19–31, 1992.

[4] R. S. Boyer and J. S. Moore. A fast string searching algorithm. *Communications of the ACM*, 20(10):762–772, 1977.

[5] M. Crochemore, A. Czumaj, L. Gasieniec, S. Jarominek, T. Lecroq, W. Plandowski, and W. Rytter. Speeding up two string-matching algorithms. *Algorithmica*, 12(4–5):247–267, 1994.

[6] I. Herms and S. Rahmann. Computing alignment seed sensitivity with probabilistic arithmetic automata. In K. Crandall and J. Lagergren, editors, *Algorithms in Bioinformatics (WABI)*, volume 5251 of *LNCS*, pages 318–329. Springer, 2008.

[7] J. Hopcroft. An $n \log n$ algorithm for minimizing the states in a finite automaton. In Z. Kohavi and A. Paz, editors, *The theory of machines and computations*, pages 189–196. Academic Press, New York, 1971.

[8] R. N. Horspool. Practical fast searching in strings. *Software-Practice and Experience*, 10:501–506, 1980.

[9] D. E. Knuth, J. Morris, and V. R. Pratt. Fast pattern matching in strings. *SIAM Journal on Computing*, 6(2):323–350, 1977.

[10] T. Knuutila. Re-describing an algorithm by Hopcroft. *Theoretical Computer Science*, 250(1-2):333–363, January 2001.

[11] G. Kucherov, L. Noé, and M. Roytberg. A unifying framework for seed sensitivity and its application to subset seeds. *Journal of Bioinformatics and Computational Biology*, 4(2):553–569, 2006.

[12] H. M. Mahmoud, R. T. Smythe, and M. Régnier. Analysis of Boyer-Moore-Horspool string-matching heuristic. *Random Structures and Algorithms*, 10(1-2):169–186, 1997.

[13] T. Marschall and S. Rahmann. Probabilistic arithmetic automata and their application to pattern matching statistics. In P. Ferragina and G. M. Landau, editors, *Combinatorial Pattern Matching (CPM)*, volume 5029 of *LNCS*, pages 95–106. Springer, 2008.

[14] T. Marschall and S. Rahmann. Efficient exact motif discovery. *Bioinformatics*, 25(12):i356–i364, 2009.

[15] T. Marschall and S. Rahmann. Exact analysis of Horspool's and Sunday's pattern matching algorithms with probabilistic arithmetic automata. In A.-H. Dediu, H. Fernau, and C. Martín-Vide, editors, *Proceedings of the 4th International Conference on Language and Automata Theory and Applications (LATA)*, volume 6031 of *LNCS*, pages 439–450, 2010.

[16] G. Navarro and M. Raffinot. *Flexible Pattern Matching in Strings*. Cambridge University Press, 2002.

[17] M. Schulz, D. Weese, T. Rausch, A. Döring, K. Reinert, and M. Vingron. Fast and adaptive variable order markov chain construction. In K. A. Crandall and J. Lagergren, editors, *Algorithms in Bioinformatics (WABI'08)*, volume 5251 of *LNCS*, pages 306–317. Springer, 2008.

[18] R. T. Smythe. The Boyer-Moore-Horspool heuristic with Markovian input. *Random Structures and Algorithms*, 18(2):153–163, 2001.

[19] D. M. Sunday. A very fast substring search algorithm. *Communications of the ACM*, 33(8):132–142, 1990.

[20] T. Tsai. Average case analysis of the Boyer-Moore algorithm. *Random Structures and Algorithms*, 28(4):481–498, 2006.

[21] S. Wu and U. Manber. A fast algorithm for multi-pattern searching. Technical report, University of Arizona, Tucson, AZ, 1994.